

# Coupling Experts and Routers in Mixture-of-Experts via an Auxiliary Loss

---

**Ang Lv**, Jin Ma, Yiyuan Ma, Siyuan Qiao

April, 2026

# Background

**Mixture-of-Experts (MoE)** uses standalone routers to select experts for each token.

MoEs suffer from *the separation* between the **router's decision-making** and the **experts' execution**.

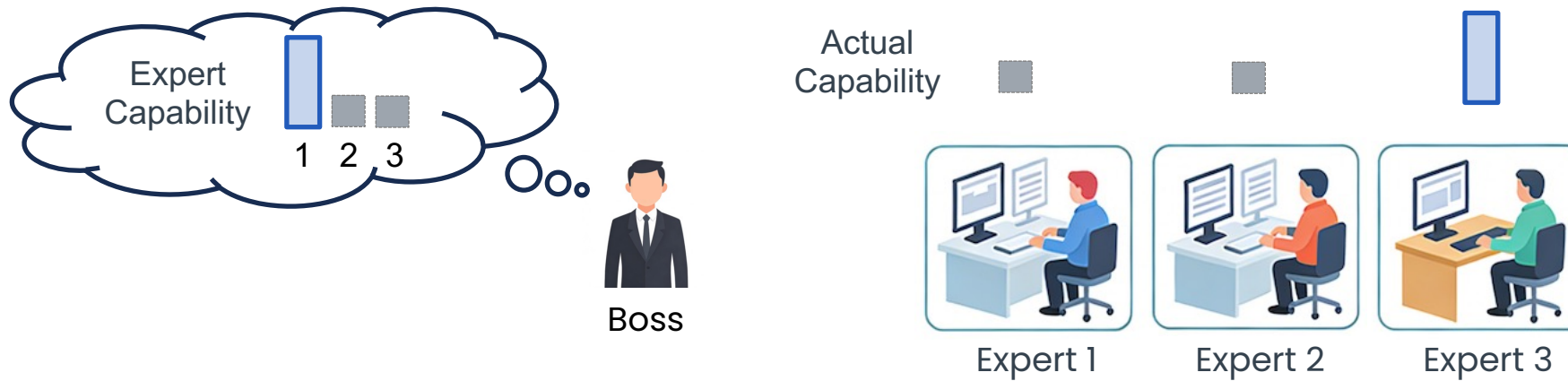


Figure: The separation between decision-making and execution.

# Background

**Mixture-of-Experts (MoE)** uses standalone routers to select experts for each token.

MoEs suffer from *the separation* between the **router's decision-making** and the **experts' execution**.

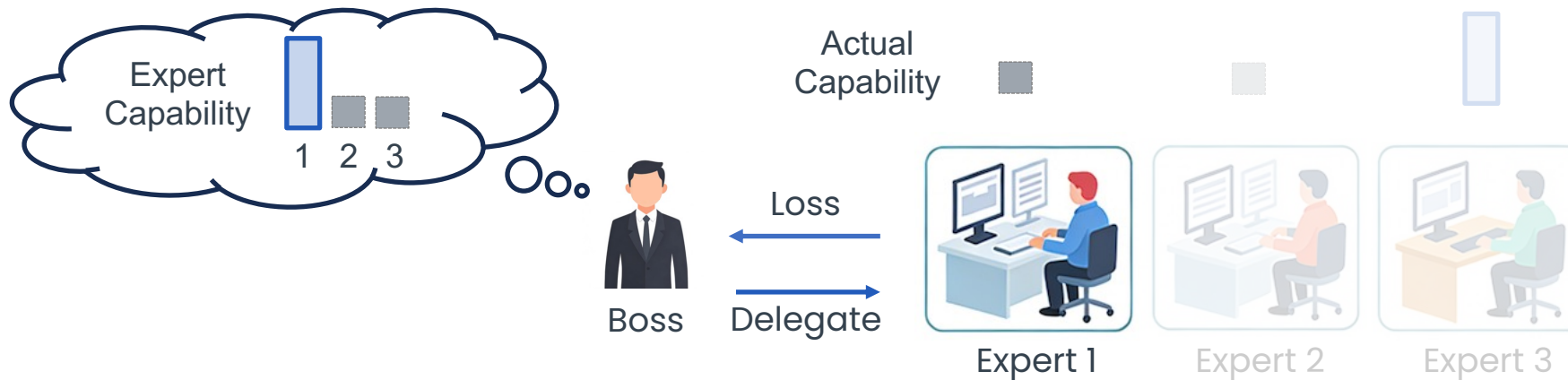


Figure: The separation between decision-making and execution.

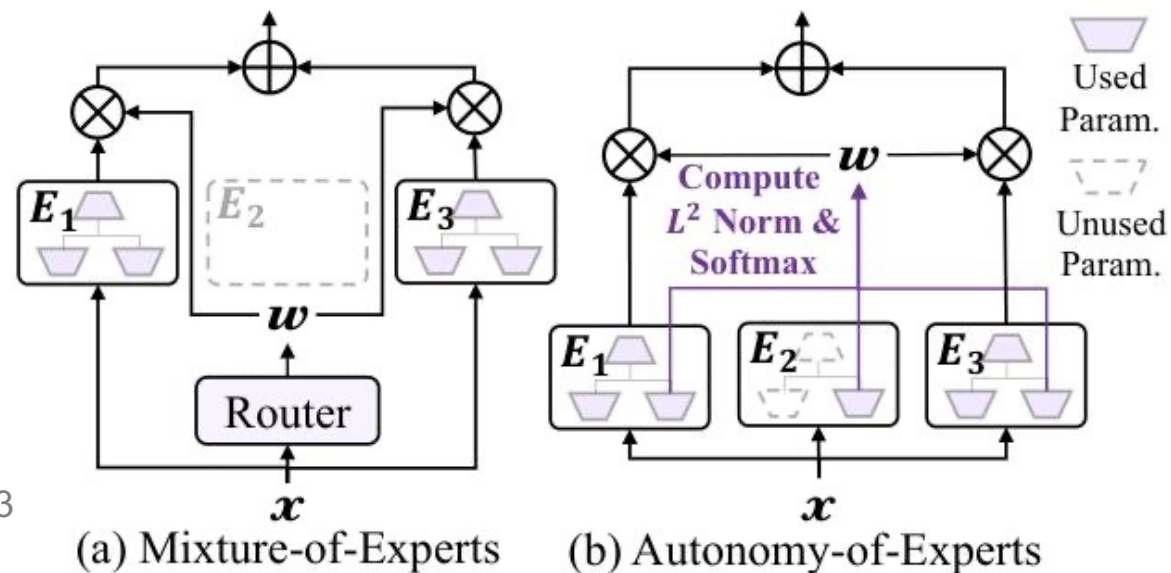
# Previous Works

**An insight:** A module matching the input's requirements yields high activations. (Section 4.1, [1])

Inspired by this, **Autonomy-of-Experts (AoE)** [2] addresses the "separation issue":

Each token is fed to all experts, and only the top-K experts with the highest intermediate activation norms are selected.

AoE is impractical due to its exhaustive computational complexity, which scales linearly with the number of tokens.



[1] Deja Vu: Contextual Sparsity for Efficient LLMs at Inference Time, ICML 2023

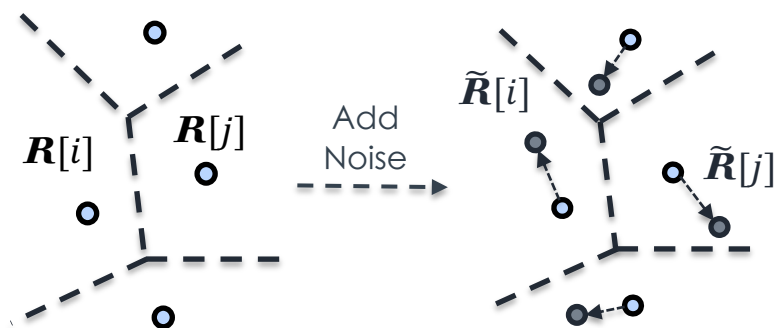
[2] Autonomy-of-Experts Models, ICML 2025

# Design Principles

1. Routers must be retained in MoE architectures to preserve routing efficiency.
2. An auxiliary loss that enables interaction between experts and routers can strengthen their coupling.
3. The loss must have complexity independent of the number of input tokens.

# Expert-Router Coupling Loss

There are three steps to calculate ERC loss, which efficiently addresses the separation issue.



Router parameters  $\mathbf{R} \in \mathbb{R}^{n \times d}$  are cluster centers.

**Step 1:** Perturb  $\mathbf{R}$  as proxy tokens:

$$\tilde{\mathbf{R}}[i] = \mathbf{R}[i] \odot \delta_i. \quad \delta_i \in \mathbb{R}^d$$

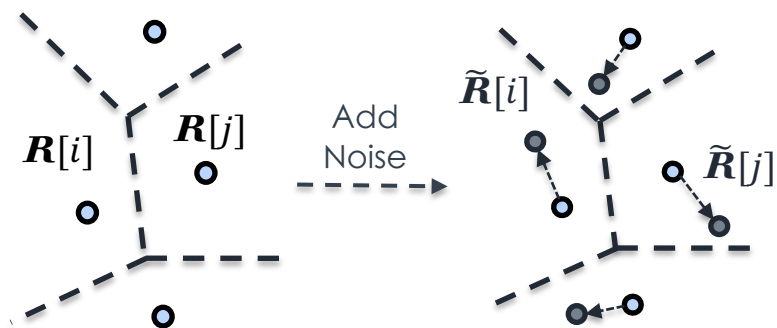
$$\delta_i \sim \mathcal{U}(1 - \epsilon_i, 1 + \epsilon_i)^d.$$

$$\epsilon_i \leq \frac{\|\mathbf{R}[i] - \mathbf{R}[j]\|}{2\|\mathbf{R}[i]\|}. \quad j = \arg \min_{j^* \neq i} \|\mathbf{R}[i] - \mathbf{R}[j^*]\|$$

$\tilde{\mathbf{R}}[i]$  is a proxy for tokens routed to expert  $i$ .

# Expert-Router Coupling Loss

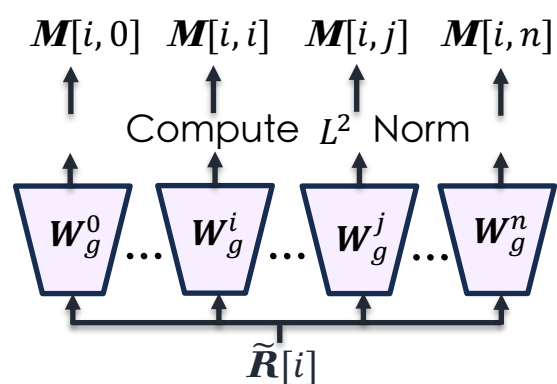
There are three steps to calculate ERC loss, which efficiently addresses the separation issue.



Router parameters are cluster centers.

**Step 1:** Perturb  $R$  as proxy tokens.

$\tilde{R}[i]$  is a proxy for tokens routed to expert  $i$ .

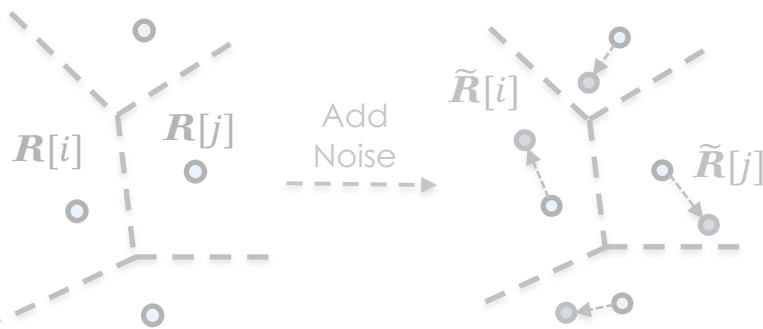


**Step 2:** Input all  $\tilde{R}[i]$  to all experts;  $W_g^i \in \mathbb{R}^{d \times D}$

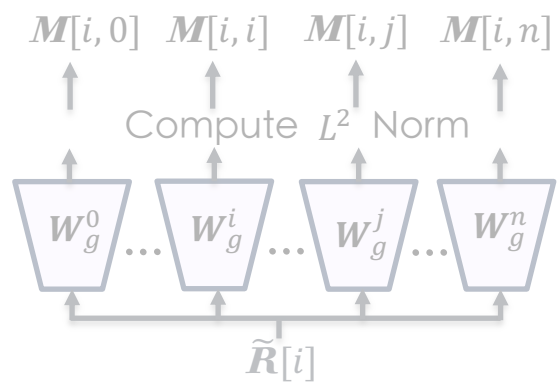
Compute  $M \in \mathbb{R}^{n \times n}$ , the  $L^2$  norm of the intermediate activations.

$M[i, j]$  is the activation norm produced by expert  $j$  given  $\tilde{R}[i]$ .

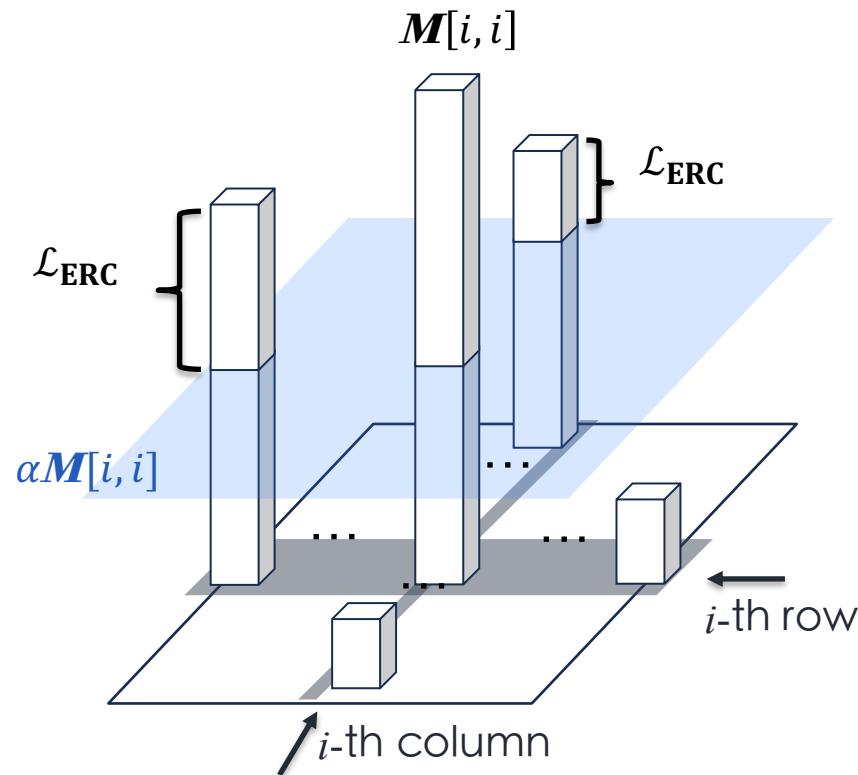
# Expert-Router Coupling Loss



Step 1: Perturb  $R$  as proxy tokens.



Step 2: Input all  $\tilde{R}[i]$  to all experts;  
Compute  $M \in \mathbb{R}^{n \times n}$ .

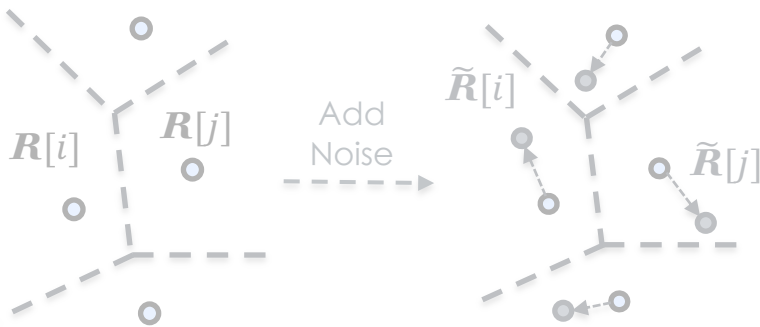


$M[i, j]$  is the activation norm produced by expert  $j$  given  $\tilde{R}[i]$ .

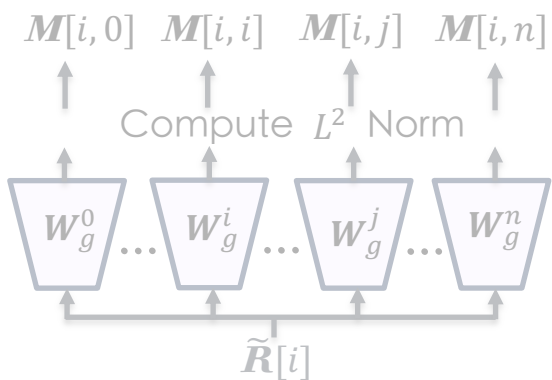
Step3: The ERC loss enforces

- (1)  $M[i, j] < \alpha M[i, i] \quad \forall i, j \neq i$
- (2)  $M[j, i] < \alpha M[i, i]$ .

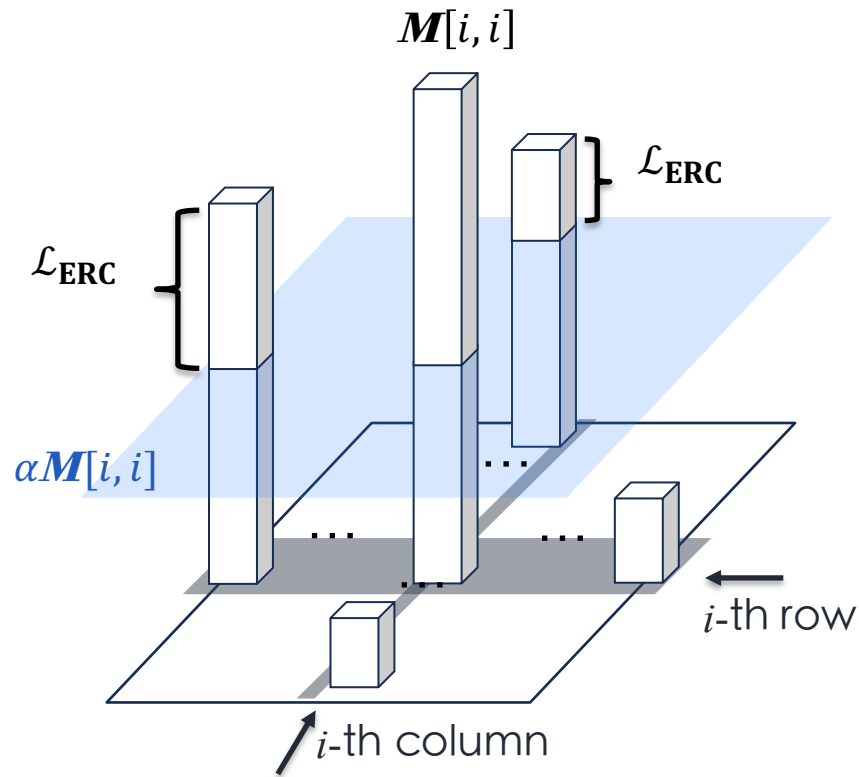
# Expert-Router Coupling Loss



Step 1: Perturb  $R$  as proxy tokens.



Step 2: Input all  $\tilde{R}[i]$  to all experts;  
Compute  $M \in \mathbb{R}^{n \times n}$ .



Step3: The ERC loss enforces

- (1)  $M[i, j] < \alpha M[i, i] \quad \forall i, j \neq i$
- (2)  $M[j, i] < \alpha M[i, i]$ .

$$\mathcal{L}_{ERC} = \frac{1}{n^2} \sum_{i=1}^n \sum_{j \neq i}^n (\max(M[i, j] - \alpha M[i, i], 0) + \max(M[j, i] - \alpha M[i, i], 0)).$$

*$n^2$  activations only, token count independent!*

# Efficiency

Additional FLOPs vs. Standard MoE,  $T$  (millions) tokens,  $K$  out of  $n$  (hundreds) experts,  $\mathbf{W}_g^i \in \mathbb{R}^{d \times D}$

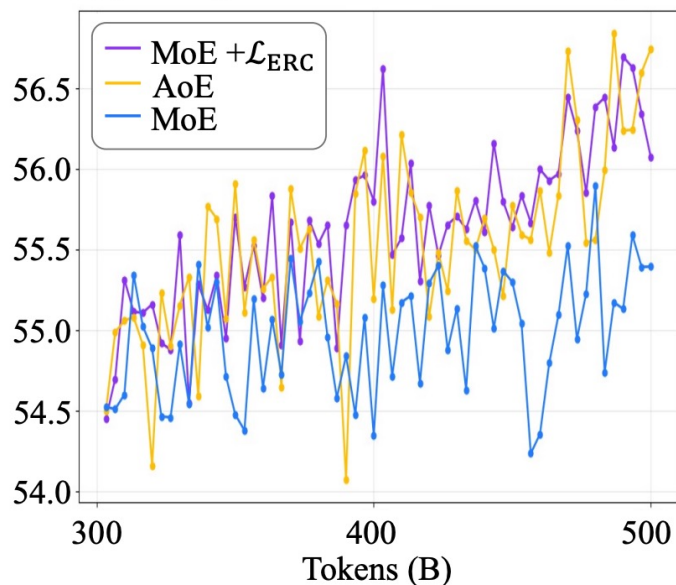
- AoE:  $2TDd(n - K)$
- ERC loss:  $2n^2Dd$ , and features **no inference overhead.**  $\mathcal{O}(Tn)$   
 $\mathcal{O}(n^2)$

Empirically, ERC loss results in a slight decrease in training throughput, with a magnitude of 0.2% to 0.8%.

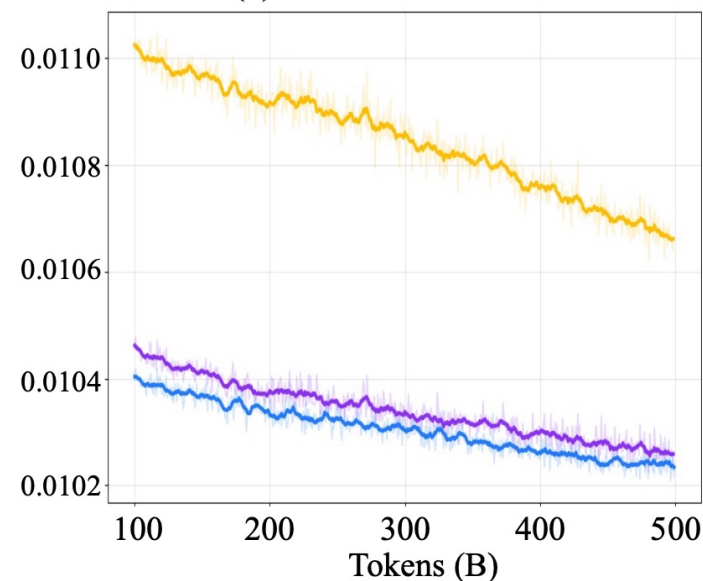
# Experiments

3B parameters (8 out of 64 experts), trained on 500B tokens, task accuracy averaged across 10 tasks.

(a) Average Downstream Task Accuracy



(b) Load Balance Loss

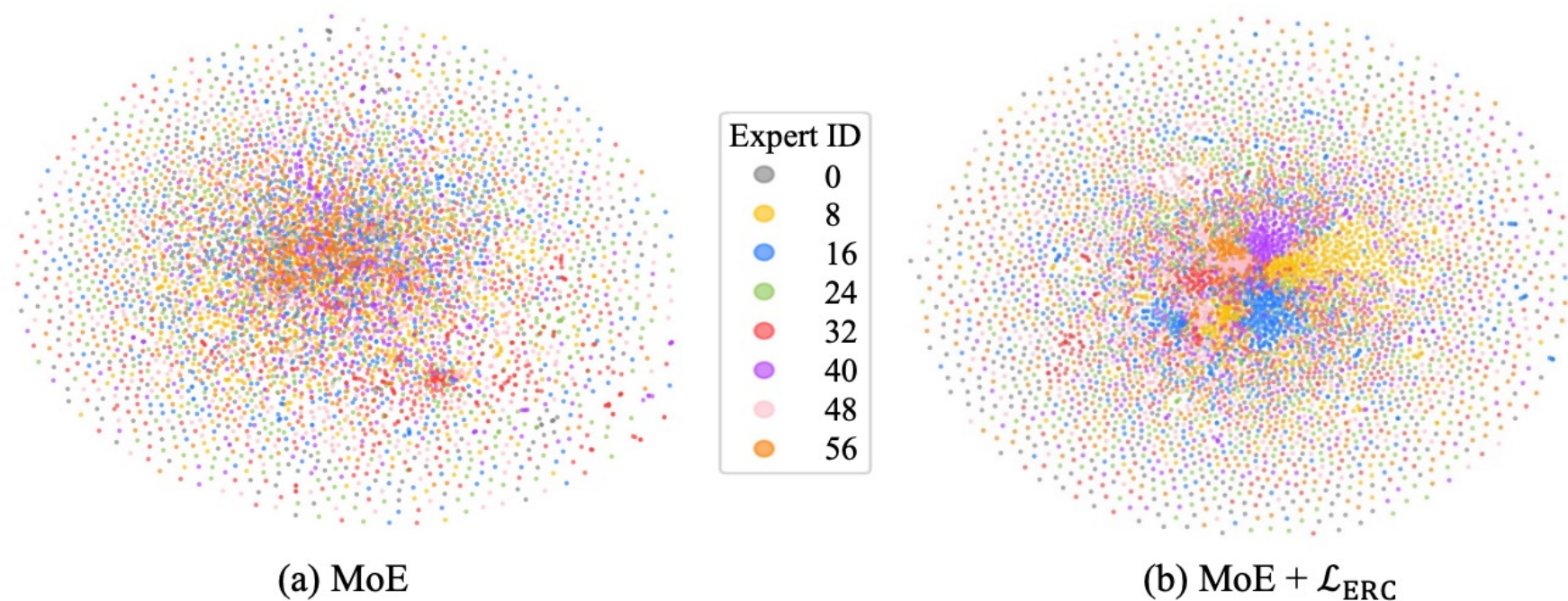


15B parameters (8 out of 256 experts), trained on 500B tokens, task accuracy averaged across 10 tasks.

	MMLU	C-Eval	MMLU-Pro	AGI-Eval	BBH	MATH	GSM8K	TriviaQA
MoE	63.2	67.5	31.0	42.0	44.3	25.7	45.2	47.2
MoE + $\mathcal{L}_{ERC}$	64.6	69.0	31.9	44.2	45.6	26.1	45.8	49.1

# ERC loss and Specialization

With the ERC loss, experts are more specialized, as they exhibit greater discrimination between tokens they process and those they do not.



**Figure 4** t-SNE projections of  $\mathbf{W}_g$  in MoE experts trained without and with the ERC loss. Our ERC loss provides greater expert specialization.

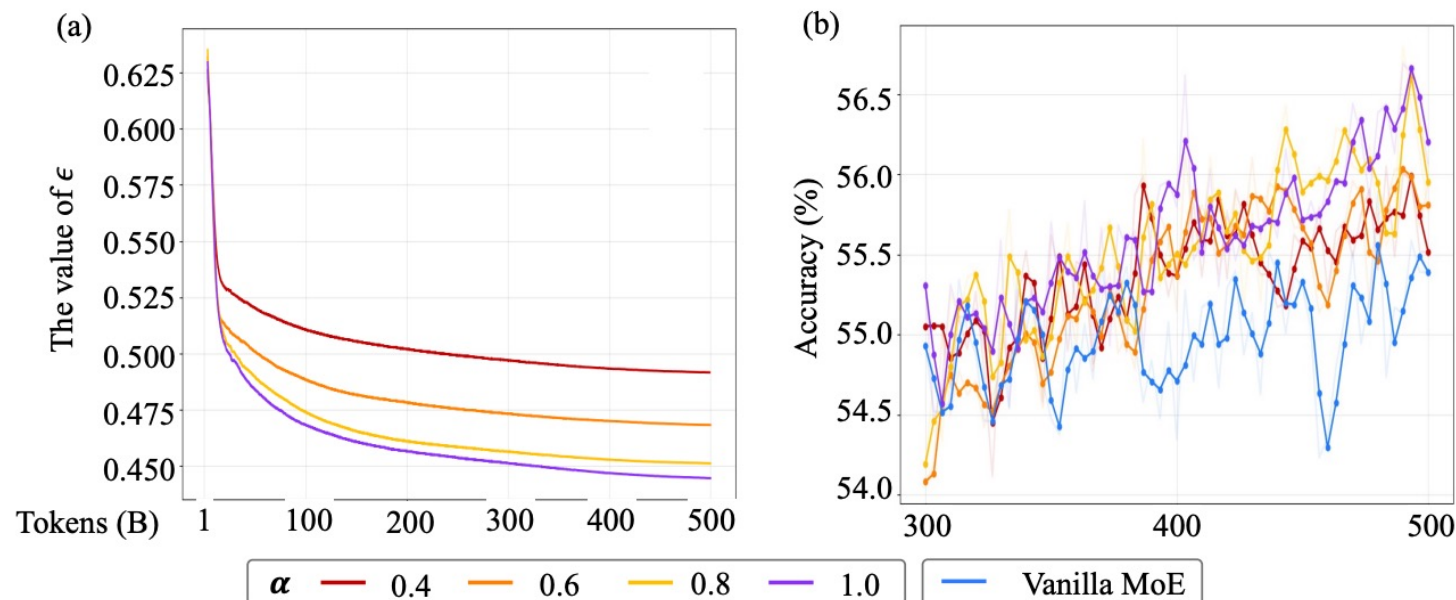
# ERC loss and Specialization

The ERC loss is an effective tool for exploring expert specialization.

(a) Since routers are deeply coupled with experts, the distance between neighboring cluster centers (i.e., the noise level  $\epsilon$ ) quantitatively reflects changes in expert specialization during training, which is controlled by  $\alpha$ .

$$\text{Recap: } \epsilon_i \leq \frac{\|\mathbf{R}[i] - \mathbf{R}[j]\|}{2\|\mathbf{R}[i]\|}, \quad j = \arg \min_{j^* \neq i} \|\mathbf{R}[i] - \mathbf{R}[j^*]\|$$

(b) Downstream performance across different values of  $\alpha$ .



# Unresolved Challenges

1. Currently, our only way to explore the optimal specialization degree is through hyperparameter tuning (e.g.,  $\alpha$ ), with performance on downstream tasks as the feedback signal.
2. Numerous factors complicate the identification of the optimal degree:
  - K and n
  - Shared expert design
  - .....
3. Is expert output orthogonality a desirable goal?

# Thank you

## Q & A

Welcome to contact me: [anglv@ruc.edu.cn](mailto:anglv@ruc.edu.cn)

Visit my homepage: [trestad.github.io](https://trestad.github.io)